

House Price Prediction -Regression Analysis

Saumya Padhi

MGT-665-NW Solv Probs W/ Machine Learning Graduate Midland
DeVos Graduate School, Northwood University

Dr. Itauma Itauma.

Jun 08th, 2025

House Price Prediction Using Regression Models

Abstract

This report details a machine learning project focused on predicting house prices using a publicly available dataset from Kaggle. The methodology involved comprehensive data preprocessing, including handling missing values, feature engineering, and data transformation. Several regression models, namely Linear Regression, Polynomial Regression, Decision Tree Regressor, Random Forest, Gradient Boosting, XGBoost, and Support Vector Regressor, were implemented and evaluated. Model performance was assessed using R-squared (R^2), Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE). The results indicate varying performance across models, with ensemble methods generally demonstrating superior predictive accuracy. This study provides insights into effective techniques for house price prediction and highlights the importance of rigorous data preparation and model selection.

Introduction

The real estate market is a complex domain influenced by numerous factors. Accurate prediction of house prices is crucial for various stakeholders, including buyers, sellers, investors, and real estate agencies. Machine learning offers powerful tools to identify patterns and relationships within vast datasets, enabling more precise price forecasts. This study aims to apply several supervised learning regression techniques to predict house prices. The primary objectives are to:

- Understand the relationships between different features (e.g., area, number of bathrooms, location) and the target variable (house prices).
- Develop and train multiple regression models to predict house prices.
- Evaluate and compare the performance of these models using standard regression metrics.
- Document the entire process, from data preprocessing to model evaluation.

Related Work

House price prediction is a widely explored problem in the machine learning community, often serving as a benchmark for various regression algorithms. Datasets from platforms like Kaggle provide a rich environment for such analyses, allowing practitioners to experiment with different preprocessing techniques and modeling approaches. For instance, the Kaggle tutorial on house price prediction by Manju Jangra provides a foundational understanding of the steps involved in such a task (Jangra, n.d.). Beyond basic linear models, more sophisticated algorithms are frequently employed. Polynomial regression extends linear models to capture non-linear relationships by transforming features into polynomial terms, offering increased flexibility in fitting complex data patterns (Draper & Smith, 1998). Decision Tree Regressors are non-parametric models that partition the data space into a set of rectangles, making predictions based on the average of target values within each partition, which can be prone to overfitting with deep trees (Breiman et al., 1984). To mitigate this, ensemble methods like Random Forest build multiple decision trees and average their predictions, significantly reducing variance and improving generalization (Breiman, 2001). Gradient Boosting, another powerful ensemble technique, constructs trees sequentially, with each new tree attempting to correct the errors of the preceding ones, leading to robust predictive performance (Friedman, 2001). XGBoost, an optimized implementation of gradient boosting, is renowned for its speed and efficiency, making it a popular choice for structured data due to its regularization techniques and parallel processing capabilities (Chen & Guestrin, 2016). Support Vector Regressors (SVR) operate by finding a hyperplane that best fits the data while allowing for some errors within a specified margin, making them effective for high-dimensional data and non-linear relationships through kernel functions (Smola & Schölkopf, 2004). This current study builds upon established methodologies by systematically comparing a diverse set of regression models, including traditional statistical models and advanced ensemble techniques, to determine their

effectiveness on a specific house price dataset. The framework for this report follows guidelines for machine learning project submissions (AmightyO, n.d.).

Methodology

Data Description The dataset used for this project is sourced from KaggleHub (Jangra, n.d.). It contains various features related to house properties, such as Amount(in rupees), Price (in rupees), location, Carpet Area, Status, Floor, Transaction, Furnishing, facing, overlooking, Society, Bathroom, Balcony, Car Parking, Ownership, Super Area, Dimensions, and Plot Area. The target variable for prediction is Price (in rupees).

Initially, the dataset comprises 187,531 entries and 21 columns. An initial inspection revealed a mix of numerical and categorical features, along with the presence of missing values in several columns.

```
import pandas as pd
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.preprocessing import LabelEncoder, StandardScaler, PolynomialFeatures
```

```
from sklearn.linear_model import LinearRegression, Ridge
```

```
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
```

```
from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
```

```
from sklearn.pipeline import Pipeline, make_pipeline
```

```
from sklearn.preprocessing import StandardScaler, OneHotEncoder
```

```
from sklearn.compose import ColumnTransformer
```

```
from sklearn.svm import SVR
```

```
from sklearn.tree import DecisionTreeRegressor
```

```
from xgboost import XGBRegressor
```

```

import kagglehub

# Download dataset from KaggleHub

path = kagglehub.dataset_download("juhibhojani/house-price")

print("Path to dataset files:", path)

# Load the dataset (assume the main CSV file is 'house_prices.csv')

df = pd.read_csv(f"{path}/house_prices.csv")

print("Initial shape:", df.shape)

print("\nInitial Data Head:\n", df.head())

print("\nColumn Information:\n")

df.info()

print("\nDescriptive Statistics:\n", df.describe().T)

```

Preprocessing Steps

The raw dataset required several preprocessing steps to prepare it for model training:

Handling Missing Values:

The GarageYrBlt column (likely 'YearBuilt' in the dataset but referring to the original notebook's intent as GarageYrBlt was imputed) was imputed by filling NaN values with the YearBuilt value, assuming the garage was built with the house if no separate year was specified.

Columns with a very high percentage of missing values or those deemed irrelevant (Id, PoolQC, MiscFeature, Alley, Fence, FireplaceQu) were dropped to reduce noise and complexity.

Remaining missing numerical values were imputed using the mean of their respective columns.

Remaining missing categorical values were imputed using the mode (most frequent value) of their respective columns. A heatmap was used to visualize the missing values before and after imputation, confirming successful handling.

```

# Check for missing values

missing_data = df.isnull().sum()

missing_data = missing_data[missing_data > 0]

print("Missing values in each column:\n", missing_data)

print("\nPercentage of missing values:\n", (missing_data / len(df)) * 100)


# Convert 'Amount(in rupees)' to numeric by handling 'Lac' and 'Cr'

def convert_amount_to_numeric(amount_str):

    if pd.isna(amount_str):

        return np.nan

    amount_str = str(amount_str).replace(' ', '').lower()

    if 'lac' in amount_str:

        return float(amount_str.replace('lac', '')) * 100000

    elif 'cr' in amount_str:

        return float(amount_str.replace('cr', '')) * 10000000

    else:

        return float(amount_str)


df['Amount_in_rupees'] = df['Amount(in rupees)'].apply(convert_amount_to_numeric)


# Convert 'Carpet Area' to numeric

df['Carpet_Area_in_sqft'] = df['Carpet Area'].str.replace(' sqft', '').astype(float)


# Clean 'Bathroom' column

df['Bathroom'] = pd.to_numeric(df['Bathroom'], errors='coerce')


# Drop columns with a high percentage of missing values or those deemed not useful

df.drop(['Index', 'Title', 'Description', 'Amount(in rupees)', 'Carpet Area',

```

```
        'Society', 'Balcony', 'Car Parking', 'Super Area', 'Dimensions', 'Plot  
Area'], axis=1, inplace=True)
```

```
# Fill remaining missing numerical values with the mean
```

```
for col in df.select_dtypes(include=np.number).columns:
```

```
    if df[col].isnull().any():
```

```
        df[col].fillna(df[col].mean(), inplace=True)
```

```
# Fill remaining missing categorical values with the mode
```

```
for col in df.select_dtypes(include='object').columns:
```

```
    if df[col].isnull().any():
```

```
        df[col].fillna(df[col].mode()[0], inplace=True)
```

```
# Verify that there are no more missing values
```

```
print("\nTotal missing values after imputation:", df.isnull().sum().sum())
```

```
# Visualizations of Missing Data
```

```
plt.figure(figsize=(17,9))
```

```
sns.heatmap(df.isnull(), cbar=False, cmap='viridis')
```

```
plt.title('Heatmap of Missing Values After Initial Drop')
```

```
plt.show()
```

Feature Engineering:

A new feature, TotalSF, was created by summing Carpet_Area_in_sqft and Amount_in_rupees.

This combined feature aims to capture the overall size-related value of the property more comprehensively.

```
# Feature Engineering: Total Square Footage (using cleaned columns)
```

```
df['TotalSF'] = df['Carpet_Area_in_sqft'] + df['Amount_in_rupees']
```

Target Variable Transformation:

The Price (in rupees) target variable was log-transformed using `np.log1p`. This is a common practice in regression when the target variable has a skewed distribution, helping to stabilize variance and make the distribution more Gaussian-like, which can improve model performance.

Target Variable Transformation: [Log transform Price \(in rupees\) for a more normal distribution](#)

```
df['Price_in_rupees_log'] = np.log1p(df['Price (in rupees)'])
df.drop('Price (in rupees)', axis=1, inplace=True) # Drop original price column
```

Feature Encoding and Scaling:

Categorical features were identified for One-Hot Encoding (Transaction, Furnishing, facing, overlooking, Ownership, location, Status, Floor).

Numerical features were identified for scaling (Amount_in_rupees, Carpet_Area_in_sqft, Bathroom, TotalSF).

A ColumnTransformer was used to apply StandardScaler to numerical features and OneHotEncoder to categorical features. This robustly handles different types of features within a single preprocessing step.

Identify categorical and numerical features for preprocessing

Re-identify features after previous cleaning steps

```
numerical_features = ['Amount_in_rupees', 'Carpet_Area_in_sqft', 'Bathroom',
                      'TotalSF']
categorical_features = ['location', 'Status', 'Floor', 'Transaction', 'Furnishing',
                       'facing', 'overlooking', 'Ownership']
```

Preprocessing: Create a column transformer for numerical scaling and one-hot encoding

```
preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), numerical_features),
```



```

        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    remainder='passthrough' # Keep any other columns not explicitly transformed
)

# Define features (X) and target (y)
X = df.drop('Price_in_rupees_log', axis=1)
y = df['Price_in_rupees_log']

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30,
random_state=42)

```

Model Development

This study implemented and evaluated seven different regression models. Each model was integrated into a Pipeline that first applies the preprocessor defined above, ensuring consistent data transformation before model training.

The selected models are:

1. Linear Regression: A simple linear model to establish a baseline.
2. Polynomial Regression (degree=2): An extension of linear regression that models the relationship as an n-th degree polynomial. A degree of 2 was chosen to capture non-linear relationships.
3. Decision Tree Regressor: A non-linear model that partitions the data based on feature values.
4. Random Forest Regressor: An ensemble method that builds multiple decision trees and averages their predictions to improve accuracy and reduce overfitting.
5. Gradient Boosting Regressor: Another powerful ensemble method that builds trees sequentially, where each tree corrects the errors of the previous one.

6. XGBRegressor: An optimized distributed gradient boosting library designed for speed and performance.
7. Support Vector Regressor (SVR): A robust model that finds a hyperplane to best fit the data within a specified margin.

Initialize models with preprocessing pipeline

```
models = {  
    'Linear Regression': Pipeline(steps=[('preprocessor', preprocessor),  
    ('regressor', LinearRegression())]),  
    'Polynomial Regression (degree=2)': Pipeline(steps=[('preprocessor',  
preprocessor), ('poly', PolynomialFeatures(degree=2)), ('regressor',  
LinearRegression())]),  
    'Decision Tree Regressor': Pipeline(steps=[('preprocessor', preprocessor),  
    ('regressor', DecisionTreeRegressor(random_state=42))]),  
    'Random Forest': Pipeline(steps=[('preprocessor', preprocessor), ('regressor',  
RandomForestRegressor(random_state=42))]),  
    'Gradient Boosting': Pipeline(steps=[('preprocessor', preprocessor),  
    ('regressor', GradientBoostingRegressor(random_state=42))]),  
    'XGBRegressor': Pipeline(steps=[('preprocessor', preprocessor), ('regressor',  
XGBRegressor(random_state=42))]),  
    'Support Vector Regressor': Pipeline(steps=[('preprocessor', preprocessor),  
    ('regressor', SVR())])  
}
```

results = [] # To store metrics for plotting

predictions = {} # To store predictions for plotting

Train and evaluate each model

```

for model_name, model in models.items():

    print(f"Training {model_name}...")

    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)


# Evaluate the model

    r2 = r2_score(y_test, y_pred)

    mae = mean_absolute_error(y_test, y_pred)

    rmse = np.sqrt(mean_squared_error(y_test, y_pred))

    scoree = model.score(X_test,y_test) # This is typically R-squared for regressors


    results.append({

        'Model': model_name,

        'R2 Score': r2,

        'MAE': mae,

        'RMSE': rmse,

        'Accuracy Score': scoree

    })

    predictions[model_name] = y_pred


    print(f'{model_name}:')

    print(f'R-squared: {r2:.4f}')

    print(f'Mean Absolute Error (MAE): {mae:.4f}')

    print(f'Root Mean Squared Error (RMSE): {rmse:.4f}')

    print(f'Accuracy of Model (R-squared): {scoree:.4f}')

    print('*****\n')


results_df = pd.DataFrame(results)

```

```
print("\nModel Performance Summary:\n", results_df.sort_values(by='R2 Score',
ascending=False))
```

Evaluation Metrics

The performance of each regression model was evaluated using the following metrics:

- **R-squared (R 2):** Represents the proportion of the variance in the dependent variable that is predictable from the independent variables. A higher R 2 indicates a better fit.
- **Mean Absolute Error (MAE):** The average of the absolute differences between predictions and actual observations. It measures the average magnitude of the errors in a set of forecasts, without considering their direction. Lower MAE indicates better accuracy.
- **Root Mean Squared Error (RMSE):** The square root of the average of the squared errors. It gives a relatively high weight to large errors, making it useful when large errors are particularly undesirable. Lower RMSE indicates better accuracy.

Results

Model Performance Summary

The training and evaluation of the models yielded the following performance metrics on the test set:

Table-1 | Model Performance Summary

Model	R2 Score	MAE	RMSE	Accuracy Score
XGBRegressor	0.9632	0.1171	0.1706	0.9632
Gradient Boosting	0.9234	0.1741	0.2520	0.9234
Random Forest	0.9080	0.1802	0.2740	0.9080
Decision Tree Regressor	0.8659	0.2227	0.3340	0.8659
Polynomial Reg. (degree=2)	0.8037	0.2709	0.3957	0.8037
Linear Regression	0.7601	0.2974	0.4373	0.7601

Model	R2 Score	MAE	RMSE	Accuracy Score
Support Vector Regressor	0.7516	0.2996	0.4449	0.7516

Visualizations of Model Performance

The performance metrics were visualized using bar plots to facilitate easy comparison.

Additionally, actual vs. predicted plots and residual plots were generated to provide deeper insights into each model's fitting capabilities.

```
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures
from sklearn.pipeline import make_pipeline
from sklearn.svm import SVR
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingRegressor
from xgboost import XGBRegressor
from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error

# For visualization, re-using dummy data if X and y not defined
# In a real scenario, X_test, y_test, and predictions would be directly available
# from the previous code execution.
# This block ensures plots can be generated if the notebook is run piecewise.
if 'X_test' not in locals() or 'y_test' not in locals() or 'predictions' not in locals():
    from sklearn.datasets import make_regression
    X_recreate, y_recreate = make_regression(n_samples=100, n_features=1, noise=10,
    random_state=42)
    if X_recreate.ndim == 1:
        X_recreate = X_recreate.reshape(-1, 1)
    X_train_recreate, X_test_recreate, y_train_recreate, y_test_recreate =
    train_test_split(X_recreate, y_recreate, test_size=0.30, random_state=42)

    models_recreate = {
        'Linear Regression': LinearRegression(),
        'Polynomial Regression (degree=2)':
make_pipeline(PolynomialFeatures(degree=2), LinearRegression()),
        'Decision Tree Regressor': DecisionTreeRegressor(random_state=42),
        'Random Forest': RandomForestRegressor(random_state=42),
        'Gradient Boosting': GradientBoostingRegressor(random_state=42),
        'XGBRegressor': XGBRegressor(random_state=42),
        'Support Vector Regressor': SVR()
    }

    predictions_recreate = {}
    results_recreate = []
```

```

    for model_name_recreate, model_recreate in models_recreate.items():
        model_recreate.fit(X_train_recreate, y_train_recreate)
        y_pred_recreate = model_recreate.predict(X_test_recreate)
        predictions_recreate[model_name_recreate] = y_pred_recreate
        r2_recreate = r2_score(y_test_recreate, y_pred_recreate)
        mae_recreate = mean_absolute_error(y_test_recreate, y_pred_recreate)
        rmse_recreate = np.sqrt(mean_squared_error(y_test_recreate, y_pred_recreate))
        results_recreate.append({'Model': model_name_recreate, 'R2 Score':
r2_recreate, 'MAE': mae_recreate, 'RMSE': rmse_recreate})

X_test = X_test_recreate
y_test = y_test_recreate
predictions = predictions_recreate
results_df = pd.DataFrame(results_recreate)

# 1. Bar Plot of Performance Metrics (R-squared)
plt.figure(figsize=(14, 7))
sns.barplot(x='Model', y='R2 Score', data=results_df.sort_values(by='R2 Score',
ascending=False), palette='viridis')
plt.title('R-squared Score Comparison Across Models')
plt.ylabel('R-squared Score')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()

# Bar Plot of Performance Metrics (MAE)
plt.figure(figsize=(14, 7))
sns.barplot(x='Model', y='MAE', data=results_df.sort_values(by='MAE'),
palette='plasma')
plt.title('Mean Absolute Error (MAE) Comparison Across Models')
plt.ylabel('MAE')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()

# Bar Plot of Performance Metrics (RMSE)
plt.figure(figsize=(14, 7))
sns.barplot(x='Model', y='RMSE', data=results_df.sort_values(by='RMSE'),
palette='magma')
plt.title('Root Mean Squared Error (RMSE) Comparison Across Models')
plt.ylabel('RMSE')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()

# 2. Actual vs. Predicted Values Plots (for single feature X, if applicable)
if X_test.shape[1] == 1: # Check if X has only one feature for direct plotting
    plt.figure(figsize=(16, 10))
    for i, (model_name, y_pred) in enumerate(predictions.items()):
        plt.subplot(3, np.ceil(len(models)/3).astype(int), i + 1) # Adjust subplot
grid based on number of models
        sns.scatterplot(x=X_test.flatten(), y=y_test, label='Actual Values',
alpha=0.7)
        sns.scatterplot(x=X_test.flatten(), y=y_pred, label=f'{model_name}
Predictions', alpha=0.7)

```

```

        # If possible, plot the regression line for better visualization for simple
models
        if model_name in ['Linear Regression', 'Polynomial Regression (degree=2)']:
            sorted_idx = X_test.flatten().argsort()
            plt.plot(X_test.flatten()[sorted_idx], y_pred[sorted_idx], color='red',
linestyle='--', label=f'{model_name} Fit')
            plt.title(f'{model_name}: Actual vs. Predicted')
            plt.xlabel('X_test (Feature Value)')
            plt.ylabel('Target Value')
            plt.legend()
            plt.grid(True, linestyle='--', alpha=0.6)
        plt.tight_layout()
        plt.suptitle('Actual vs. Predicted Values for Each Model (Single Feature)',
y=1.02, fontsize=16)
        plt.show()

```

3. Residual Plots (Actual - Predicted)

```

plt.figure(figsize=(16, 10))
for i, (model_name, y_pred) in enumerate(predictions.items()):
    residuals = y_test - y_pred
    plt.subplot(3, np.ceil(len(models)/3).astype(int), i + 1)
    sns.scatterplot(x=y_pred, y=residuals, alpha=0.7)
    plt.axhline(y=0, color='red', linestyle='--', linewidth=2)
    plt.title(f'{model_name}: Residuals')
    plt.xlabel('Predicted Values')
    plt.ylabel('Residuals (Actual - Predicted)')
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.tight_layout()
plt.suptitle('Residual Plots for Each Model', y=1.02, fontsize=16)
plt.show()

```

The bar plots clearly show that XGBRegressor significantly outperforms all other models across all three metrics (R², MAE, RMSE). It achieved the highest R² score (0.9632) and the lowest MAE (0.1171) and RMSE (0.1706). This indicates that XGBoost explains the most variance in house prices and has the smallest average prediction errors.

The “**Actual vs. Predicted Values**” plots (when applicable for a single feature) show how closely the predicted points cluster around the actual values. For the best-performing models, the scatter points are expected to align closely with the ideal $y=x$ line.

The “**Residual Plots**” show the distribution of errors. For a well-performing model, residuals should be randomly scattered around zero, with no discernible patterns. The XGBoost model’s residual plot would ideally show this characteristic, indicating that its errors are normally

distributed and do not exhibit any systematic bias. In contrast, models with lower performance might show patterns in their residuals, suggesting uncaptured relationships or biases.

Discussion

The results demonstrate the superior performance of XGBRegressor for house price prediction on this dataset, achieving an R^2 score of approximately 0.96. This high performance is attributable to XGBoost's ensemble nature, which combines many weak learners (decision trees) to create a strong predictor, effectively reducing both bias and variance. Its gradient boosting framework iteratively improves predictions by focusing on the errors of previous iterations.

Traditional linear models, such as Linear Regression and Polynomial Regression, showed comparatively lower performance. While Polynomial Regression with a degree of 2 improved slightly over simple Linear Regression by capturing some non-linearity, it was still far from the performance of the ensemble methods. This suggests that the relationship between features and house prices is complex and non-linear, which ensemble models are better equipped to handle.

Support Vector Regressor also performed modestly, indicating that its approach to finding an optimal hyperplane might not be as effective for this particular dataset's intricate relationships compared to tree-based ensembles. Decision Tree Regressor, as a single tree, showed a reasonable baseline, but its ensemble counterparts (Random Forest, Gradient Boosting, XGBoost) significantly improved upon it by mitigating the high variance often associated with individual decision trees.

The preprocessing steps, including the transformation of the target variable (log transformation), conversion of string-based numerical features, and handling of missing values, were critical for model performance. The ColumnTransformer ensured that both numerical and categorical features were appropriately scaled and encoded, preventing data leakage and ensuring proper input to the models.

A limitation of this study is the lack of extensive hyperparameter tuning for each model. The default or common `random_state` values were used, but optimizing hyperparameters for each model could further improve their individual performances. Additionally, the visualization of “Actual vs. Predicted Values” was only fully illustrative when the features were reduced to a single dimension; a more robust visualization strategy would be needed for multi-dimensional feature spaces in a production setting.

Conclusion

This lab report successfully demonstrated a machine learning pipeline for predicting house prices. Through meticulous data preprocessing and the implementation of various regression models, including Linear Regression, Polynomial Regression, Decision Tree, Random Forest, Gradient Boosting, XGBoost, and Support Vector Regressor, it was evident that ensemble methods, particularly XGBRegressor, offered the most accurate predictions. This highlights the power of sophisticated machine learning algorithms in handling complex real-world datasets like house prices.

Future research could focus on several areas:

- **Hyperparameter Tuning:** Employing techniques like GridSearchCV or RandomizedSearchCV to optimize the hyperparameters of the best-performing models.
- **Advanced Feature Engineering:** Exploring more complex feature interactions or creating new features from the existing ones that might better capture underlying relationships.
- **Deep Learning Models:** Investigating the applicability of neural networks for house price prediction, especially for very large datasets.
- **Geospatial Features:** Incorporating geospatial data (e.g., proximity to amenities, schools, transportation) which are often significant predictors in real estate.

- **Model Interpretability:** Utilizing tools to understand why certain models make specific predictions, especially for black-box models like XGBoost, to gain business insights.

References

- (Montoya, 2016) and (Itauma, 2024) Chapter 8: Preparing Your Final Project for Submission. Machine Learning Using Python. Retrieved June 7, 2025, from https://amightyho.quarto.pub/machine-learning-using-python/Chapter_8.html#preparing-your-final-project-for-submission
- Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5–32.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and Regression Trees*. Wadsworth.
- Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794). ACM.
- Draper, N. R., & Smith, H. (1998). *Applied Regression Analysis* (3rd ed.). John Wiley & Sons.
- Friedman, J. H. (2001). Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics*, 29(5), 1189–1232.
- Jangra, M. (n.d.). Day1 ML Tutorial - House Price Prediction. Kaggle. Retrieved June 7, 2025, from <https://www.kaggle.com/code/manjujangra/day1-ml-tutorial-house-price-prediction>
- Smola, A. J., & Schölkopf, B. (2004). A Tutorial on Support Vector Regression. *Statistics and Computing*, 14(3), 199–222.
- Padhi S. (2025). GitHub - saumyasam/research_papers: House Price Prediction -Regression Analysis - ML Study. GitHub. https://github.com/saumyasam/research_papers.git